



# Exact computation of the Maximum Entropy Potential of spiking neural networks models

Rodrigo Cofre, Bruno Cessac

## ► To cite this version:

Rodrigo Cofre, Bruno Cessac. Exact computation of the Maximum Entropy Potential of spiking neural networks models. 2014. hal-00861397v4

**HAL Id: hal-00861397**

**<https://hal.inria.fr/hal-00861397v4>**

Preprint submitted on 14 May 2014

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Exact computation of the Maximum Entropy Potential of spiking neural networks models

R. Cofré<sup>1</sup> and B. Cessac<sup>1</sup>

<sup>1</sup>*NeuroMathComp team (INRIA, UNSA LJAD) 2004 Route des Lucioles, 06902 Sophia-Antipolis, France*

(Dated: April 8, 2014)

Understanding how stimuli and synaptic connectivity influence the statistics of spike patterns in neural networks is a central question in computational neuroscience. Maximum Entropy approach has been successfully used to characterize the statistical response of simultaneously recorded spiking neurons responding to stimuli. But, in spite of good performance in terms of prediction, the fitting parameters do not explain the underlying mechanistic causes of the observed correlations. On the other hand, mathematical models of spiking neurons (neuro-mimetic models) provide a probabilistic mapping between stimulus, network architecture and spike patterns in terms of conditional probabilities. In this paper we build an exact analytical mapping between neuro-mimetic and Maximum Entropy models.

PACS numbers: 87.19.lo 05.10.-a 87.10.-e 87.85.dq

## I. INTRODUCTION

The spike response of a neural network to external stimuli is largely conditioned by the stimulus itself, synaptic interactions and neuronal network history [17]. Understanding these dependences is a current challenge in neuroscience [26]. Since spikes occur irregularly both within and over repeated trials [27], it is reasonable to characterize spike trains using statistical methods and probabilistic descriptions.

Among existing approaches to characterize spike train statistics, the Maximum Entropy Principle (MaxEnt) [21] has been successfully applied to data from the cortex and the retina [14, 24, 28, 29]. It consists of fixing a set of constraints, determined as the empirical average of observables measured from spiking activity. Maximizing the statistical entropy, given those constraints, provides a unique probability called Gibbs distribution. The choice of constraints determines a “model”. Prominent examples are the Ising model [28, 29] where constraints are firing rates and probabilities that 2 neurons fire at the same time, the Ganmor-Schneidman-Segev model [14], which considers additionally the probability of triplets and quadruplets of spikes, or the Tkačik et al model [31] where the probability that  $K$  out of  $N$  cells in the network generate simultaneous action potentials is an additional constraint. In these examples the statistics has no memory and successive times are considered statistically independent; but Markovian models where the probability of a spike pattern at a given time depends on the spike history can be considered as well [24, 32]. MaxEnt models depends on a set of parameters (Lagrange multipliers) which are fitting parameters. In statistical physics language, these are parameters conjugated to the constraints, just like the inverse temperature  $\beta = \frac{1}{kT}$  is conjugated to the energy, or chemical potential is conjugated with the number of particles. However, whereas inverse temperature or chemical potential have a clear in-

terpretation thanks to the links between thermodynamics and statistical physics, the fitting parameters (Lagrange multipliers) used for spike train statistics do not benefit from such deep relations and are interpreted e.g. via loose analogies with magnetic systems. For example the Lagrange multipliers  $J_{ij}$  conjugated to pairwise spike coincidence are interpreted as “functional interactions” [14] due to their analogy with magnetic interactions in the Ising model. Likewise the parameters  $h_i$  conjugated with single spike events (whose average is the firing rate) are believed to be related with an effective stimulus received by neuron  $i$ . However, the connection between “functional” interactions  $J_{ij}$  and real interactions (e.g. synapses) in the network remains elusive, as well as the link between effective stimuli and the stimulus viewed by a neuron.

An alternative approach is based on spiking neuron models, providing a mathematical description of neural dynamics. These models give a probabilistic mapping between network architecture, stimuli, spiking history of the network and spiking response in terms of conditional probabilities of spike pattern given the network history. Prominent examples are the Linear-Non Linear model (LN), the Generalized Linear Model (GLM) [6, 11] or Integrate-and-Fire models [17]. In all these examples the conditional probabilities that a spike pattern occurs at time  $t$  given the network spike history are explicit functions of “structural” parameters in the neural network (that can be interpreted as synaptic weights  $\mathcal{W}$  matrix, and stimulus  $\mathcal{I}$ ) (Fig. 1a).

These conditional probabilities define a Markov process that mimics the biophysical dynamics of neurons in a network and the mechanisms that govern spike trains emission, including stimulus dependence and neurons interactions via synapses. We call them *neuro-mimetic* statistical models.

To summarize, at least two different representations can be used to analyze spike train statistics in neural

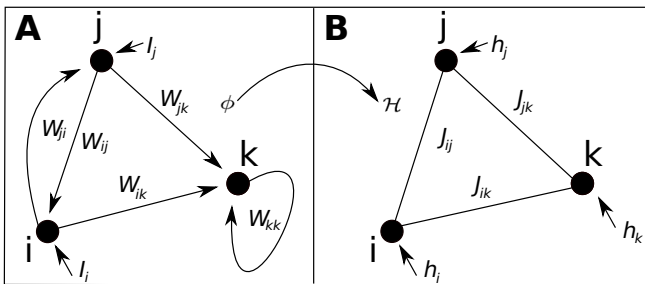


FIG. 1: (A) Neuro-mimetic approach. Neurons are interacting via synaptic weights  $W_{ij}$  and submitted to a stimulus  $I$ . Spike probabilities are explicit functions of these parameters. (B) MaxEnt statistical approach. Here the relation between neurons are expressed by functional parameters allowing to correctly fit the correlations in the model. The graph represents the Ising model where only local fields and pairwise interactions are drawn. More general interactions are considered in the text. In Ising model pairwise interactions are symmetric (represented without arrows). We are looking for an explicit and exact correspondence between these two representations.

networks (fig. 1). The goal of this paper is to establish an explicit and exact correspondence between these two representations.

A previous result attempting to describe such a relationship can be found in [12]. Here, the authors fit a leaky Integrate and Fire model matching spike train data from a population of retinal ganglion cells. At the same time they fit a MaxEnt Ising model from this data. This allows them to compare in particular synaptic weights  $W_{ij}$  with MaxEnt Ising couplings  $J_{ij}$ . Another work in this direction can be found in [18] in which stimulus dependent MaxEnt is introduced based on (LN) model, attempting to include stimulus information into the “local fields” of the Ising model. Both examples are limited to the Ising model, thus do not include memory effects in the MaxEnt statistics.

We propose here a generalization which allows us to handle more general types of neuro-mimetic models as well as general spatio-temporal MaxEnt distributions (including memory). The method we used is based on Hammersley-Clifford decomposition [20] and periodic orbit invariance from ergodic theory [25]. The techniques are therefore different from [12, 18].

More generally, we answer the following questions:

**Question 1:** *Given an ergodic Markov process, where the transition probabilities are known, can we construct a MaxEnt potential, with a minimum of constraints, reproducing exactly the (spatio-temporal) statistics of this process?*

This is the most general question we answer in this paper. It is important to notice that our results are not

restricted to spike trains and neural networks, but to any ergodic Markov chain. The next question focuses on the correspondence between MaxEnt parameters and structural parameters defining spiking neural networks.

**Question 2:** *Given the transition probabilities from a neuro-mimetic model, is it possible to derive an analytic correspondence between MaxEnt fitting parameters (Lagrange multipliers) and neuro-mimetic structural parameters? For example can we establish a correspondence between the local fields  $h_i$  and the external stimulus  $I_i$ ? Between Ising couplings  $J_{ij}$  with  $W_{ij}$ , the synaptic weights?*

We show that there exists an exact and analytic correspondence revealing that Lagrange multipliers are complex non-linear functions of structural parameters. For example the local field  $h_i$  or the functional interactions  $J_{ij}$  are non-linear functions depending *generically* on all stimuli and all synaptic weights.

Additionally this correspondence raises up the question of dimensionality. A neuro-mimetic model with  $N$  neurons typically has  $N^2 + N$  independent parameters ( $N^2$  synaptic weights and  $N$  stimuli); a MaxEnt model with memory depth  $D$  may have up to  $2^{N(D+1)}$  independent parameters (see text). The dimensionality of these two types of models is drastically different. When mapping a MaxEnt model to a neuro-mimetic model there is clearly a loss of dimensionality.

**Question 3:** *Consider a MaxEnt model equivalent to a neuro-mimetic model. Then the difference in dimensionality between the spaces of parameters of both models suggests that either many of MaxEnt parameters are zero, or that there are relations among them, i.e. they are not independent. What is the generic situation?*

The paper is organized as follows. In section II we introduce some notations and present MaxEnt models with spatio-temporal constraints. The introduction of memory requires to define the Gibbs distribution in a more general setting than usual statistical physics definition. Although this formalism is well known [16], it allows us to make the connection between Gibbs distributions and Markov chains, a necessary step in our construction. In section III we develop our method, based on equivalence between potentials, Hammersley-Clifford hierarchy and periodic orbits invariance. In section IV we present an example based on a discrete time Integrate and Fire model in which we compute explicitly the “local fields”  $h_i$  and “Ising couplings”  $J_{ij}$  as non linear functions of  $W, I$ . We finally present the conclusions of this work in which we address especially the issue of the difference of dimensionality between the space of parameters of MaxEnt and neuro-mimetic models.

## II. SETTING

We study a network composed by  $N$  neurons. Time has been discretized so that a neuron can at most fire a spike within one time bin. A spike train is represented by a binary time series with entries  $\omega_k(n) = 1$  if neuron  $k$  fires at time  $n$  and  $\omega_k(n) = 0$  otherwise. The *spiking pattern* at time  $n$  is the vector  $\omega(n) = [\omega_k(n)]_{k=1}^N$ . A *spike block*  $\omega_{n_1}^{n_2}$  is an ordered list of spiking patterns where spike times range from  $n_1$  to  $n_2$ .

A *potential*  $\mathcal{H}$  of range  $R = D + 1$  is a function that associates to each spike block  $\omega_0^D$  a real value. We assume  $\mathcal{H}(\omega_0^D) > -\infty$ . Any such potential can be written:

$$\mathcal{H}(\omega_0^D) = \sum_{l=0}^L h_l m_l(\omega_0^D), \quad (1)$$

where  $L = 2^{NR} - 1$ . The  $h_l$ 's are real numbers whereas the function  $m_l$  with  $m_l(\omega_0^D) = \prod_{u=1}^r \omega_{k_u}(n_u)$  is called a *monomial*.  $m_l(\omega_0^D) = 1$  if and only if neurons  $k_u$  fires at time  $n_u$  and zero otherwise.  $r$  is the *degree* of the monomial. By analogy with spin systems, we see from (1) that monomials somewhat constitute formal *spatio-temporal* interactions: degree 1 monomials corresponds to “self-interactions”, degree 2 to pairwise interactions, and so on. The  $h_l$ 's characterize the intensity of the corresponding interaction.

In many examples most  $h_l$ 's are equal to zero. For instance, Ising model considers only monomials of the form  $\omega_i(0)$  (singlets) or  $\omega_i(0)\omega_j(0)$  (instantaneous pairwise events).

$$\mathcal{H}(\omega_0^D) = \sum_{i=0}^N h_i \omega_i(0) + \sum_{i=0}^N \sum_{j>i}^N J_{ij} \omega_i(0) \omega_j(0). \quad (2)$$

More generally, the potential (1) considers spatio-temporal events occurring within a time horizon  $R$ . This potential defines a unique stationary probability  $\mu$ , called *equilibrium state with potential*  $\mathcal{H}$  satisfying the following variational principle [16, 22]:

$$\mathcal{P}[\mathcal{H}] = \sup_{\nu \in \mathcal{M}} (\mathcal{S}[\nu] + \nu[\mathcal{H}]) = \mathcal{S}[\mu] + \mu[\mathcal{H}], \quad (3)$$

where  $\mathcal{M}$  is the set of stationary probabilities defined on the set of spike trains, whereas:

$$\mathcal{S}[\nu] = - \lim_{n \rightarrow \infty} \frac{1}{n+1} \sum_{\omega_0^n} \nu[\omega_0^n] \log \nu[\omega_0^n],$$

is the entropy of the probability  $\nu$ . The average of  $\mathcal{H}$  with respect to  $\nu$  is noted  $\nu[\mathcal{H}] = \sum_{l=0}^L h_l \nu[m_l]$ . In (3) we denote  $\mu(\mathcal{H})$ ,  $(\nu(\mathcal{H}))$  the average of  $\mathcal{H}$  with respect to  $\mu$ ,  $(\nu)$ . We insist on a point: we only consider stationary (time translation invariant) probabilities in this paper.

The quantity  $\mathcal{P}[\mathcal{H}]$  is called *free energy* and has the following properties:

- $\mathcal{P}[\mathcal{H}]$  is a log generating function of cumulants. We have:

$$\frac{\partial \mathcal{P}[\mathcal{H}]}{\partial h_l} = \mu[m_l],$$

the average of  $m_l$  with respect to  $\mu$  and:

$$\frac{\partial^2 \mathcal{P}[\mathcal{H}]}{\partial h_k \partial h_l} = \frac{\partial \mu[m_l]}{\partial h_k} = \sum_{n=-\infty}^{+\infty} C_{m_k, m_l}(n), \quad (4)$$

where  $C_{m_k, m_l}(n)$  is the correlation function between the two monomials  $m_k$  and  $m_l$  at time  $n$  in the equilibrium state  $\mu$ . Note that correlation functions decay exponentially fast whenever  $\mathcal{H}$  has finite range and  $\mathcal{H} > -\infty$ , thus  $\sum_{n=-\infty}^{+\infty} C_{m_k, m_l}(n) < +\infty$ . Eq. (4) characterizes the variation in the average value of  $m_l$  when varying  $h_k$  (linear response). The corresponding matrix is a susceptibility matrix. It controls the Gaussian fluctuations of monomials around their mean (central limit theorem) [5]. When considering potential of range 1 ( $D = 0$ ) eq (4) reduces to the classical fluctuation-dissipation theorem, because the corresponding process has no memory (successive times are independent thus  $C_{m_k, m_l}(n) = 0$  unless  $n = 0$ ).

- $\mathcal{P}(\mathcal{H})$  is a convex function of  $h_l$ 's. This ensures the uniqueness of the solution of (3).

### Transfer Matrix

We now show that any potential  $\mathcal{H} > -\infty$  with the form (1) is naturally associated with a Markov chain whose invariant distribution is the equilibrium state  $\mu$  satisfying the variational principle (3).  $\mu$  is additionally a Gibbs distribution.

We first recall that a Markov chain is defined by a set of transition probabilities  $P[\omega(n) | \omega_{n-D}^{n-1}]$ . We assume here that the memory depth of the chain  $D$  is constant and finite, although an extension of the present formalism to variable length Markov chains ( $D$  variable) or chains with complete connections ( $D$  infinite) is possible [13]. We also assume that the chain is homogeneous (transition probabilities do not depend explicitly on time) and primitive (there exist  $n$  such that any two states are connected by a path of length  $n$ , with positive probability) [33]. Then the Markov chain admits a unique invariant probability  $\mu$  which obeys the Chapman-Kolmogorov relation:  $\forall n_1, n_2, n_2 + D - 1 > n_1$ ,

$$\mu[\omega_{n_1}^{n_2}] = \prod_{n=n_1}^{n_2-D} P[\omega(n+D) | \omega_n^{n+D-1}] \mu[\omega_{n_1}^{n_1+D-1}]. \quad (5)$$

Introducing:

$$\phi(\omega_n^{n+D}) = \log P[\omega(n+D) | \omega_n^{n+D-1}], \quad (6)$$

we have

$$\mu[\omega_{n_1}^{n_2}] = e^{\sum_{n=n_1}^{n_2-D} \phi(\omega_n^{n+D})} \mu[\omega_{n_1+D-1}^{n_1+D}].$$

$\phi$  is called a *normalized potential*. To each  $\mathcal{H}$  of the form (1) corresponds a unique normalized potential  $\phi$  and a unique invariant measure  $\mu$ . Although this correspondence can be found in many text-books (see for example [16, 22]), we summarize it here since it is the core of our approach.

### From $\mathcal{H}$ to Markov chains

Each spike block is associated to a unique integer (index)  $l = \sum_{k=1}^N \sum_{n=0}^D 2^{nN+k-1} \omega_k(n)$ , where neurons  $k = 1, \dots, N$  are considered from top to bottom and time  $n = 0, \dots, D$  from left to right in the spike train. We denote  $\omega^{(l)}$  the spike block corresponding to the index  $l$ . Here an example with  $N = 2$  and  $R = 3$ ,  $\omega^{(6)} = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix}$ .

Consider two spike blocks  $\omega^{(l)}, \omega^{(l')}$  of range  $D \geq 1$ . The transition  $\omega^{(l)} \rightarrow \omega^{(l')}$  is *legal* if  $\omega^{(l)}, \omega^{(l')}$  have a common block  $\omega_1^{D-1}$ . Here is an example of a legal transition:

$$\omega^{(l)} = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix}; \omega^{(l')} = \begin{bmatrix} 0 & 1 & 1 \\ 0 & 1 & 0 \end{bmatrix}.$$

and a forbidden transition:

$$\omega^{(l)} = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix}; \omega^{(l')} = \begin{bmatrix} 0 & 1 & 1 \\ 0 & 1 & 0 \end{bmatrix}.$$

Any block  $\omega_0^D$  of range  $R = D + 1$  can be viewed as a legal transition from the block  $\omega^{(l)} = \omega_0^{D-1}$  to the block  $\omega^{(l')} = \omega_1^D$  and in this case we write  $\omega_0^D \sim \omega^{(l)} \omega^{(l')}$ .

We construct the *transfer matrix*  $\mathcal{L}$ :

$$\mathcal{L}_{\omega^{(l)}, \omega^{(l')}} = \begin{cases} e^{\mathcal{H}(\omega_0^D)} & \text{if } \omega_0^D \sim \omega^{(l)} \omega^{(l')} \\ 0, & \text{otherwise.} \end{cases}$$

This is a  $2^{ND} \times 2^{ND}$  matrix whose indexes are spike blocks. Note that, from the assumption  $\mathcal{H} > -\infty$ , each legal transition corresponds to a positive entry in the matrix  $\mathcal{L}$ . Therefore  $\mathcal{L}$  is primitive and satisfies the Perron-Frobenius theorem [15].

As a consequence of the Perron-Frobenius theorem,  $\mathcal{L}$  has a unique real positive eigenvalue  $s$ , strictly larger in modulus than the other eigenvalues, and with right,  $R$ , and left,  $L$ , eigenvectors:  $\mathcal{L}R = sR$ ,  $L\mathcal{L} = sL$ . The following holds:

(a) These eigenvectors have strictly positive entries  $R(\cdot) > 0$ ,  $L(\cdot) > 0$ ; their arguments are blocks of range  $D$ . They can be chosen so that the scalar product  $\langle L, R \rangle = 1$ .

(b) The following potential:

$$\phi(\omega_0^D) = \mathcal{H}(\omega_0^D) - \log R(\omega_0^{D-1}) + \log R(\omega_1^D) - \log s \quad (7)$$

is normalized i.e. it defines via (6) an homogeneous Markov chain with transition probability  $P[\omega(D) | \omega_0^{D-1}] = e^{\phi(\omega_0^D)}$ .

(c) The unique invariant probability of this Markov chain is:

$$\mu(\omega_0^{D-1}) = R(\omega_0^{D-1}) L(\omega_0^{D-1}). \quad (8)$$

(d) It follows from Chapman-Kolmogorov equation (5) and from (7,8) that, for  $D > 0$ :

$$\mu[\omega_0^n] = \frac{e^{\sum_{k=0}^{n-D} \mathcal{H}(\omega_k^{k+D})}}{s^{n-D+1}} R(\omega_{n-D+1}^n) L(\omega_0^{D-1}). \quad (9)$$

(e)  $\mu$  obeys the variational principle (3) and

$$\mathcal{P}[\mathcal{H}] = \log s.$$

When considering a normalized potential, the transfer matrix becomes a stochastic transition matrix with maximal eigenvalue 1. Thus  $\mathcal{P}[\phi] = 0$ .

(f) It follows from (9) that

$\exists A, B > 0$  such that, for any block  $\omega_0^n$  the Gibbs distribution reads [5, 22]:

$$A \leq \frac{\mu[\omega_0^n]}{e^{-(n-D+1)\mathcal{P}(\mathcal{H})} e^{-\sum_{k=0}^{n-D} \mathcal{H}(\omega_k^{k+D})}} \leq B. \quad (10)$$

This is actually the definition of Gibbs distributions in ergodic theory [34].

This definition encompasses the classical definition of Gibbs distributions,  $\frac{e^{\mathcal{H}}}{Z}$  found in standard textbooks of statistical physics. Let us indeed consider a potential of range  $R = 1$ , ( $D = 0$ ). This is a limit case in the definition of the transfer matrix where transitions between spike patterns  $\omega(0) \rightarrow \omega(1)$  are considered and where all transitions are legal.  $\mathcal{L}_{\omega(0), \omega(1)} = e^{\mathcal{H}(\omega(0))}$ , thus each row has the form:

$$(e^{\mathcal{H}(\omega(0))}, e^{\mathcal{H}(\omega(0))}, \dots, e^{\mathcal{H}(\omega(0))}).$$

The matrix  $\mathcal{L}$  is degenerated with a maximum eigenvalue:

$$s = Z = \sum_{\omega(0)} e^{\mathcal{H}(\omega(0))}$$

and all other eigenvalues 0. The left eigenvector corresponding to  $s = Z$  is:

$$L = \left( \frac{1}{Z}, \frac{1}{Z}, \dots, \frac{1}{Z} \right)$$

whereas  $R(\omega(0)) = e^{\mathcal{H}(\omega(0))}$ . Note that we have normalized  $L$  so that  $\langle L, R \rangle = 1$ . We have therefore  $\mu(\omega(0)) = \frac{e^{\mathcal{H}(\omega(0))}}{Z}$ , the classical form for the Gibbs distribution. The normalized potential in the limiting case is  $\phi(\omega(0)) = \mathcal{H}(\omega(0)) - \log Z$ , whereas the Markov chain has no memory: successive events are independent. This last remark reflects a central weakness of memory-less MaxEnt models to describe neuron dynamics. They neither involve memory nor time causality. In order to consider realistic situations, where a spike pattern probability depends on the past spike history, MaxEnt models must be constructed as we did [35]. Indeed, it is not possible to extend the form  $\frac{e^{\mathcal{H}(\omega(0))}}{Z}$  to spatio-temporal potentials. This is obvious from equation (9):  $\mu[\omega_0^n]$  has not the form  $\frac{e^{\mathcal{H}(\omega(0))}}{Z_n}$  with  $Z_n = \sum_{\omega_0^n} e^{\mathcal{H}(\omega_0^n)}$ . This can also be readily seen from (10). However the following holds:

$$\mathcal{P}[\mathcal{H}] = \lim_{n \rightarrow \infty} \frac{1}{n} \log Z_n.$$

This outlines a crucial point: as soon as one introduces memory in the MaxEnt, infinite time limits have to be considered in order to fully characterize the statistics. This is *mutatis mutandis* the same procedure as taking the thermodynamic limit in spatial lattices [16].

### Equivalent potentials

Although a potential  $\mathcal{H}$  of the form (1) corresponds (if  $\mathcal{H} > -\infty$ ) to a unique normalized potential  $\phi$  and Gibbs distribution  $\mu$ , this correspondence is not one to one. To a normalized potential  $\phi$  corresponds infinitely many potentials of the form (1). Hence two potentials  $\mathcal{H}^{(1)}, \mathcal{H}^{(2)}$  can correspond to the same Gibbs distribution (We call them equivalent).

A standard result in ergodic theory states that  $\mathcal{H}^{(1)}$  and  $\mathcal{H}^{(2)}$  are equivalent if and only if there exists a range  $D > 0$  function  $f$  such that [5]:

$$\mathcal{H}^{(2)}(\omega_0^D) = \mathcal{H}^{(1)}(\omega_0^D) - f(\omega_0^{D-1}) + f(\omega_1^D) + \Delta, \quad (11)$$

where  $\Delta = \mathcal{P}[\mathcal{H}^{(2)}] - \mathcal{P}[\mathcal{H}^{(1)}]$ . This relation establishes a *strict* equivalence and does not correspond e.g. to renormalization. The validity of (11) can be readily seen by plugging  $\mathcal{H}^{(2)}$  in the variational formula (3) the terms corresponding to  $f$  cancels because  $\nu$  is time-translation invariant. Therefore, the supremum is reached for the same Gibbs distribution as  $\mathcal{H}^{(1)}$  whereas  $\Delta$  is indeed the

difference of free energies. The “only if” part is more tricky.

Equation (7) is a particular case of equation (11), where  $\mathcal{H}^{(2)} = \phi, \mathcal{H}^{(1)} = \mathcal{H}, f = \log R$  and  $\Delta = -\log s$ . This equation has the virtue to unify two very different approaches. It establishes a relation between Markov chain normalized potentials (6) on one hand and potentials of the form (1) on the other hand (the arrow  $\phi \rightarrow \mathcal{H}$  in fig. 1). Equation (11) answers therefore the first part of the question (1), but, by itself does not provide a straightforward way to exploit it, due to the arbitrariness in the choice of  $f$ . Indeed, there are infinitely many potentials  $\mathcal{H}$  corresponding to the same Gibbs distribution (the same normalized potential  $\phi$ ).

This arbitrariness in the choice of  $f$  raises a natural question closely related to the second part of question (1). Given a normalized potential is it possible to find, among the infinite family of equivalent potentials, a canonical form of  $\mathcal{H}$  with a minimal number of terms? The situation is a bit like normal forms in bifurcations theory where variable changes allows one to eliminate locally non resonant terms in the Taylor expansion of the vector field [2]. Here, the role of the variable changes is played by  $f$ . By suitable choices of  $f$  one should be able to eliminate some monomials in the expansion (1). An evident situation corresponds to monomials related by time translation, e.g.  $\omega_i(0)$  and  $\omega_i(1)$ : since any  $\nu \in \mathcal{M}$  is time translation invariant  $\nu[\omega_i(0)] = \nu[\omega_i(1)]$ , the firing rate of neuron  $i$ . Such monomials correspond to the same constraint in (3) and can therefore be eliminated. A potential where monomials, related by time translation, have been eliminated (the corresponding  $h_l$  vanishes) is called *canonical*. A canonical potential contains thus, in general,  $2^{NR} - 2^{N(R-1)}$  terms. We now show that canonical potentials cannot be further reduced.

### Canonical interactions cannot be eliminated using the equivalence equation (11)

Assume that we are given two potentials  $\mathcal{H}^{(1)}, \mathcal{H}^{(2)}$  in the canonical form, where  $\mathcal{H}^{(1)}$  has a zero coefficient for the canonical interaction  $m_l$  whereas  $\mathcal{H}^{(2)} = \mathcal{H}^{(1)} + h_l m_l$ ,  $h_l \neq 0$ . Let us show that these two potentials are not equivalent. For this we need to introduce a bit of notations further used in the text.

Since a monomial is defined by a set of spike events  $(k_u, n_u)$ , one can associate to each monomial a spike block or “mask” where the only bits ‘1’ are located at  $(k_u, n_u)$ ,  $u = 1, \dots, r$ . This mask has therefore an index. Whereas the labeling of monomials in (1) was arbitrary,  $m_l$  denotes from now on the monomial with mask  $\omega^{(l)}$ . Let us define the block *inclusion*  $\sqsubseteq$ , by  $\omega_0^D \sqsubseteq \omega_1^D$  if  $\omega_k(n) = 1 \Rightarrow \omega'_k(n) = 1$ , with the convention that the block of degree 0,  $\omega^{(0)}$ , is included in all blocks. Then, for two integers  $l, l'$ :

## Hammersley-Clifford hierarchy

$$m_l(\omega^{(l)}) = 1 \text{ if and only if } \omega^{(l')} \sqsubseteq \omega^{(l)}. \quad (12)$$

Now, from (11),  $\mathcal{H}^{(2)} = \mathcal{H}^{(1)} + h_l m_l$  and  $\mathcal{H}^{(1)}$  are therefore equivalent if one can find a  $D$ -dimensional function  $f$  such that,  $\forall \omega_0^D$ :

$$f(\omega_1^D) - f(\omega_0^{D-1}) + \Delta + h_l 1_{\omega^{(l)} \sqsubseteq \omega_0^D} = 0,$$

where  $1_{\omega^{(l)} \sqsubseteq \omega_0^D}$  is the standard indicator function that takes value 1 when  $\omega^{(l)} \sqsubseteq \omega_0^D$  and 0 otherwise. Let us consider 2 specific blocks. The block only composed by '1's contains all other blocks, and it is translation invariant so that the terms involving  $f$  cancel in the equation above. We have therefore  $\Delta + h_l = 0$ . The block only composed by '0's is also translation invariant and, if  $l > 0$  we obtain  $\Delta = 0$ , so that  $h_l = 0$ , in contradiction with the hypothesis. Therefore, two canonical potentials are equivalent *if and only if* all their canonical coefficients  $h_l$ ,  $l > 0$ , are equal [9].

There is still an arbitrariness due to the term  $h_0$  ("Gauge" invariance). One can set it equal 0 without loss of generality. In this way, there is only one canonical potential, with a minimal number of monomials, corresponding to a given stationary Markov chain.

In the next section we show how to compute the coefficients of the canonical potential  $\mathcal{H}^{(2)}$  equivalent to a known  $\mathcal{H}^{(1)}$ .

### III. METHOD

Given a spike block  $\omega^{(l_0)}$ , a *periodic orbit* of period  $\kappa$  is a sequence of spike blocks  $\omega^{(l_n)}$  where  $\omega^{(l_{p\kappa+n})} = \omega^{(l_n)}$ ,  $p \geq 0$ ,  $0 \leq n \leq \kappa - 1$ . From equation (11) we have, for such a periodic orbit,

$$\sum_{n=0}^{\kappa-1} \mathcal{H}^{(2)}(\omega^{(l_n)}) = \sum_{n=0}^{\kappa-1} \mathcal{H}^{(1)}(\omega^{(l_n)}) + \kappa\Delta, \quad (13)$$

because the  $f$ -terms disappear when summed along a periodic orbit. It follows that the sum of a potential along a periodic orbit is an invariant (up to the constant term  $\kappa\Delta$ ) in the class of equivalent potentials. This is a classical result in ergodic theory extending to infinite range potentials [23]. This equation is valid whatever periodic orbit is considered. It is singularly useful if one takes advantage of an existing hierarchy between blocks and between monomials, the Hammersley-Clifford (H-C) hierarchy [20] that we explain now.

The construction of our method is based on the (H-C) factorization theorem, proved in the seminal although unpublished paper [20]. Later simpler proofs were given in [3, 19]. This result establishes the equivalence between Markov random fields and Gibbs distributions. It was proved in the context of undirected graphs where the clique structures provide the factorization of the potential. Our result is based on a decomposition of the potential over inclusions  $\sqsubseteq$  of (spatio-temporal) blocks defined previously. Inclusions provide a hierarchical structure similar to the blackening algebra of (H-C). However (H-C) theorem does not provide by itself an explicit method to obtain from a Markov chain the corresponding canonical MaxEnt potentials. On the opposite, our method of periodic orbits allows to perform this computation.

We can express  $\mathcal{H}^{(2)}$  in the form (1), then using (12) it follows that (13) becomes:

$$\sum_{n=0}^{\kappa-1} \sum_{l \sqsubseteq l_n} h_l^{(2)} m_l(\omega^{(l_n)}) = \sum_{n=0}^{\kappa-1} \mathcal{H}^{(1)}(\omega^{(l_n)}) + \kappa\Delta \quad (14)$$

where, with a slight abuse of notations  $l \sqsubseteq l_n$  stands for  $\omega^{(l)} \sqsubseteq \omega^{(l_n)}$ .

The interesting fact about this representation is that the l.h.s of this equation is written entirely in terms of blocks included in the blocks considered in the periodic orbits. Therefore in order to compute all the coefficients  $h_l$ 's that characterize the canonical MaxEnt potential we can proceed by first obtaining the coefficient of degree 0, then the coefficients of degree 1, 2, and so on. We use equation (14) to compute from a known  $\mathcal{H}^{(1)}$  potential its associated canonical potential  $\mathcal{H}^{(2)}$ . From now on we focus in the particular case when  $\mathcal{H}^1 = \phi$  is a normalized potential. To alleviate notation we note  $\mathcal{H}^{(2)} = \mathcal{H}$ .

#### Degree 0: Free energy

Start from the first mask in hierarchy, the mask  $\omega^{(0)}$  containing only 0's, whose corresponding monomial is  $m_0 = 1$  and consider its periodic orbit, of period  $\kappa = 1$ ,  $\{\omega^{(0)}\}$ . The application of (14) gives  $h_0 = \phi(\omega^{(0)}) + \mathcal{P}[\mathcal{H}]$  and since we choose  $h_0 = 0$  for the canonical potential we obtain a direct way to compute the free energy of  $\mathcal{H}$ .

$$\mathcal{P}[\mathcal{H}] = -\phi(\omega^{(0)}) \quad (15)$$

#### Degree 1: Local Fields

Let us now consider masks of degree 1:

$$\omega^{(l_0)} = \begin{bmatrix} 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix},$$

(where dots correspond to 0) corresponding to the monomial  $\omega_i(D)$ . Also consider the periodic orbit obtained by a  $R$ -circular shift of this block ( $\kappa = R$ ):

$$\omega^{(l_0)} = \begin{bmatrix} 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix} \rightarrow \omega^{(l_1)} = \begin{bmatrix} 0 & \dots & 0 & 0 \\ \vdots & \ddots & \vdots & \vdots \\ 0 & \dots & 1 & 0 \\ \vdots & \ddots & \vdots & \vdots \\ 0 & \dots & 0 & 0 \end{bmatrix}$$

$$\dots \rightarrow \omega^{(l_D)} = \begin{bmatrix} 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix}.$$

Since the corresponding monomials of the blocks in the orbit  $\omega^{(l_0)}, \omega^{(l_1)}, \dots, \omega^{(l_D)}$  are related by time translation they correspond to the same constraint in (3). The coefficient of all but one of these monomials is therefore set to 0 in the canonical potential  $\mathcal{H}$ . We use the convention to keep the monomial  $m_{l_0}$  whose mask contains a 1 in the right most column. This convention extends to the monomials considered below. The block  $\omega^{(l)}$  considered to generate this periodic orbit has one spike corresponding to neuron  $i$ . To make this explicit we note  $h_l \equiv \mathbf{h}_i$ : Then, applying (14) to this periodic orbit we obtain:

$$\mathbf{h}_i = \phi(\omega^{(l_0)}) + \phi(\omega^{(l_1)}) + \dots + \phi(\omega^{(l_D)}) + R\phi(\omega^{(0)}) \quad (16)$$

We have thus obtained the coefficient corresponding to the monomial  $\omega_i(0)$  which is precisely  $\mathbf{h}_i$  in Ising model (2).

Considering a different  $\omega^{(l_0)}$  of degree 1 and the periodic orbit generated by its  $R$ -circular shift we get another local field term. We do the same for the  $N$  neurons.

## Degree 2: Instantaneous pairwise interactions

Let us now consider instantaneous pairwise interactions. We consider masks of the form :

$$\omega^{(l_0)} = \begin{bmatrix} 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix},$$

corresponding to the monomial  $\omega_i(D)\omega_j(D)$ , the procedure is the same as above i.e. take the periodic orbit:

$$\omega^{(l_0)} = \begin{bmatrix} 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix} \rightarrow \omega^{(l_1)} = \begin{bmatrix} 0 & \dots & 0 & 0 \\ \vdots & \ddots & \vdots & \vdots \\ 0 & \dots & 1 & 0 \\ \vdots & \ddots & \vdots & \vdots \\ 0 & \dots & 0 & 0 \end{bmatrix}$$

$$\dots \rightarrow \omega^{(l_D)} = \begin{bmatrix} 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix}.$$

The coefficients corresponding to this monomials are  $J_{ij}$  in the Ising model (2). We have, from (14):

$$J_{ij} = \sum_{n=0}^{R-1} \phi(\omega^{(\sigma^n l)}) + R\phi(\omega^{(0)}) - \sum_{n=0}^{R-1} \sum_{l'_n \subset \sigma^n l} h_{l'_n}. \quad (17)$$

For blocks  $l'_n \subset \sigma^n l$  of degree 1 the spike is either on neuron  $i$  or neuron  $j$ . The contribution of these blocks is  $\mathbf{h}_i + \mathbf{h}_j$ . In the blocks  $l'_n \subset \sigma^n l$  there is also the block  $\omega^{(0)}$ , whose contribution is  $h_0 = 0$ . Therefore, we finally have:

$$J_{ij} = \sum_{n=0}^{R-1} \phi(\omega^{(\sigma^n l)}) + R\phi(\omega^{(0)}) - \mathbf{h}_i - \mathbf{h}_j. \quad (18)$$

## Degree 2: (1 time-step memory):

For the one step of memory pairwise coefficients (e.g.  $\omega_j(0)\omega_i(1)$ ) the situation is slightly different;

$$\omega^{(l_0)} = \begin{bmatrix} 0 & \dots & 0 & 0 \\ \vdots & \ddots & \vdots & \vdots \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & 0 \end{bmatrix},$$

Here the periodic orbit generated by the  $R$ -circular shift of  $\omega^{(l_0)}$  is:

$$\omega^{(l_0)} = \begin{bmatrix} 0 & \dots & 0 & 0 \\ \vdots & \ddots & \vdots & \vdots \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & 0 \end{bmatrix} \rightarrow \omega^{(l_1)} = \begin{bmatrix} 0 & \dots & 0 & 0 \\ \vdots & \ddots & \vdots & \vdots \\ 0 & \dots & 0 & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & 0 \end{bmatrix}$$



$$\omega^{(l_2)} = \begin{bmatrix} 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 1 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 0 & 1 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 0 & 0 \end{bmatrix} \rightarrow \dots \rightarrow \omega^{(l_D)} = \begin{bmatrix} 0 & 0 & 0 & \dots \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 1 & 0 & \dots \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 1 & \dots \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 0 & \dots \end{bmatrix}$$

This orbit is not sufficient because it contains 2 unknowns in eq (14), namely the first and second blocks  $\omega^{(l_0)}$  and  $\omega^{(l_1)}$  correspond to monomials  $\omega_j(D-1)\omega_i(D)$  and  $\omega_j(D)\omega_i(0)$  which are not related by time translation, so correspond to different canonical constraints both having degree 2. Therefore it is not possible to solve (14) just generating one circular periodic orbit. Fortunately is possible to circumvent this problem by generating additional periodic orbits.

Let us now consider the following periodic orbit:

$$\begin{aligned} \omega^{(l_0)} &= \begin{bmatrix} 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 0 & 0 \end{bmatrix} \rightarrow \omega^{(l_1)} = \begin{bmatrix} 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 0 & 0 \end{bmatrix} \\ \omega^{(l_2)} &= \begin{bmatrix} 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 0 & 0 \end{bmatrix} \rightarrow \omega^{(l_3)} = \begin{bmatrix} 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 0 & 1 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 0 & 0 \end{bmatrix} \\ \omega^{(l_4)} &= \begin{bmatrix} 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 1 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 0 & 1 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 0 & 0 \end{bmatrix} \rightarrow \dots \rightarrow \omega^{(l_{R+2})} = \begin{bmatrix} 0 & 0 & 0 & \dots \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 1 & 0 & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 1 & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots \end{bmatrix} \quad (19) \end{aligned}$$

We have generated a periodic orbit, where (14) has only one unknown, the coefficient associated to the first block. All the other blocks in the orbit are either of lower degree, thus we have already computed them; or are time translations of the first block, thus their coefficient is set to zero.

This is a particular example of a general procedure that we describe now. It allows to compute hierarchically any  $h_l$ . The procedure is general, but we illustrate it with:

$$\omega^{(l)} = \begin{bmatrix} 0 & 0 & 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 1 \end{bmatrix}$$

- Step 1. Shift circularly  $\omega^{(l)}$  until the left-most spiking pattern has at least a 1. Each of the circular

shifts generate a mask, which corresponds to the same constraint in (3) so the corresponding  $h_l$  coefficient is set to zero.

$$\begin{bmatrix} 0 & 0 & 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 0 & 1 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 1 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 \end{bmatrix}$$

- Step 2. Continue circularly left shifting but, before shifting, remove the 1 with the lower neuron index, on the left most spike pattern. Tag the 1s that has been removed. Do this until the total number of left shifts including step 1 and 2 is  $R$ .

$$\begin{bmatrix} 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 0 \end{bmatrix} \\ \rightarrow \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 \end{bmatrix}$$

- Step 3. Same as step 1. All the masks generated at this step correspond to the same constraint and thus have a zero coefficient.

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 \end{bmatrix}$$

- Step 4. Do the opposite of what was done in step 2: Restore the 1's that has been removed on the left most spike pattern while circularly shifting. In this way we finally regenerate  $\omega^{(l)}$ .

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 \end{bmatrix} \\ \rightarrow \begin{bmatrix} 0 & 0 & 0 & 1 & 1 & 0 \\ 1 & 0 & 0 & 1 & 0 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 0 & 0 & 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 1 \end{bmatrix}$$

As claimed we have generated a periodic orbit where all monomials, but  $\omega^{(l)}$ , have either a coefficient 0 or have a degree smaller than  $\omega^{(l)}$  and have therefore been already computed. Obviously, when getting to larger and larger degrees the method becomes rapidly intractable because of the exponential increase in the number of terms. The hope is that the influence of monomials decays rapidly with their degree. Additionally, applying it to real data where transition probabilities are not exactly known leads to severe difficulties. These aspects will be treated in a separated paper. Our goal here was to answer the first question asked in the introduction. This goal is now achieved.

We now switch to the second question.

### From neuro-mimetic models to normalized potentials

In neuro-mimetic models the probability that the spike pattern  $\omega(n)$  occurs at time  $n$  is the result of the complex membrane potentials dynamics [17]. A simplification consists of assuming that this probability is only

a function of the spike history up to a certain memory depth  $D$ . In this framework it is possible to consider a Markov chain in which the set of states consists of spike blocks  $\omega_0^{D-1}$  from which legal transitions between blocks provides a family of conditional probabilities  $P_n[\omega(n) | \omega_{n-D}^{n-1}]$  that may depend explicitly on time as indicated by the sub-index  $n$ . In neuro-mimetic models these probabilities depends on parameters that mimics biophysical quantities such as synaptic weights matrix  $\mathcal{W}$  and stimulus  $\mathcal{I}$ . A particularly important example is the Generalized Linear Model (GLM), which assimilates the spike response as an inhomogeneous point process, with “conditional intensity”  $\lambda_k(t|H_t)$  which modulates the probability that neuron  $k$  emits a spike between times  $t$  and  $t + dt$  given  $H_t$  the history of spikes up to time  $t$ . This function is given by [1]:

$$\lambda_k(t|H_t) = f \left( b + K \cdot i(t) + \sum_j M_{kj} \cdot \omega_j \right), \quad (20)$$

where  $f$  is a non linear function;  $b$  is a vector fixing the baseline firing rate of neurons;  $K$  is a causal, time-translation invariant, linear convolution kernel that mimics a linear receptive field of neurons;  $i(t)$  is a stimulus;  $M_{kj}$  is a memory kernel that describes excitatory or inhibitory post spike effects of pre-synaptic neurons  $j$  on post-synaptic neuron  $k$ .  $\omega_j$  is the spike train of neuron  $j$ . Considering a discretization of time (time-steps  $\Delta t$ ) and a memory cut-off (memory  $D \rightarrow H_{n-D}^{n-1}$ ) from the conditional intensity we can get the conditional probability that neuron  $k$  fire at time  $n$ :

$$P[\omega_k(n) = 1 | H_{n-D}^{n-1}] \approx \lambda_k(n | H_{n-D}^{n-1}) \Delta t := p_k(n).$$

A crucial assumption in (GLM) is the conditional independence between neurons given the history  $H_t$ . Note that our formalism does not make this assumption and could be extended to neuro-mimetic models violating this assumptions (e.g models with gap junctions [10]). Here for simplicity we stick at models with conditional independence. As a consequence of the conditional independence the probability of a spike pattern reads:

$$P[\omega(n) | \omega_{n-D}^{n-1}] \approx \prod_{k=1}^N p_k(n)^{\omega_k(n)} (1 - p_k(n))^{1 - \omega_k(n)},$$

and taking logarithm we get the normalized potential (6):

$$\begin{aligned} \phi(\omega_{n-D}^n) = & \sum_{k=1}^N \left[ \omega_k(n) \log p_k(n) \right. \\ & \left. + (1 - \omega_k(n)) \log(1 - p_k(n)) \right], \end{aligned} \quad (21)$$

We can use equation (14) with  $\phi = \mathcal{H}^{(1)}$  to compute from  $\phi$  its canonical potential  $\mathcal{H}^{(2)} = \mathcal{H}$ :

$$\sum_{n=0}^{\kappa-1} \sum_{l \subseteq l_n} h_l m_l(\omega^{(l_n)}) = \sum_{n=0}^{\kappa-1} \phi(\omega^{(l_n)}) + \kappa \mathcal{P}[\mathcal{H}]. \quad (22)$$

Using this equation with the appropriate periodic orbits considered in the hierarchical order we obtain the corresponding canonical potential  $\mathcal{H}$  from the normalized potential characterizing the GLM (21). We show in the following section using a different example (discrete time Leaky Integrate and Fire model) how the coefficients corresponding to firing rates, instantaneous and 1-step memory pairwise correlations can be calculated explicitly in terms of synaptic weights and stimulus.

#### IV. THE DISCRETE TIME LEAKY INTEGRATE AND FIRE MODEL

In this section we illustrate our result in a stochastic leaky Integrate-and-Fire model with noise and stimulus [30] analyzed rigorously in [7].

This model is a discretization of the usual leaky Integrate-and-Fire model. Its dynamics reads:

$$V(t+1) = F(V(t)) + \sigma_B B(t), \quad (23)$$

where  $V(t) = (V_i(t))_{i=1}^N$  is the vector of neuron’s membrane potential at time  $t$ ;  $F(V)$  is a vector-valued function with entries:

$$F_i(V) = \gamma V_i (1 - S[V_i]) + \sum_{j=1}^N W_{ij} S[V_j] + I_i, \quad i = 1 \dots N$$

where  $\gamma \in [0, 1]$ , is the (discrete-time) “leak rate [36]”;  $S$  is a function characterizing the neuron’s firing: for a firing threshold  $\theta > 0$ ,  $S(x) = 1$  whenever  $x \geq \theta$  and  $S(x) = 0$  otherwise;  $I_i$  is an external current. In the most general version of this model,  $I_i$  depends on time. Here, we focus on the case where  $I_i$  is constant, ensuring the stationarity of dynamics. Finally, in (23),  $\sigma_B > 0$  is a variable controlling the noise intensity, where the vector  $B(t) = (B_i(t))_{i=1}^N$  is an additive noise. It has Gaussian independent and identically distributed entries with zero mean and variance 1.

#### The normalized potential

The normalized potential of the model (23) has infinite range. Indeed, a neuron has memory only back to the last time when it has fired. But this time is unbounded (although the probability that the last firing time arises before time  $m$  decays exponentially fast as  $m \rightarrow -\infty$ ).

Nevertheless, the exact potential can be approximated by the finite range potential [7].

$$\phi(\omega_0^D) = \sum_{k=1}^N \left[ \omega_k(D) \log \pi(X_k(\omega_0^{D-1})) + (1 - \omega_k(D)) \log(1 - \pi(X_k(\omega_0^{D-1}))) \right], \quad (24)$$

where the function  $\pi$  is:

$$\pi(x) = \frac{1}{\sqrt{2\pi}} \int_x^{+\infty} e^{-\frac{u^2}{2}} du.$$

$\phi$  has therefore the form of a (GLM) (21). All functions appearing below depend on the spike block  $\omega_0^{D-1}$  and make explicit the dependence of the network state (membrane potentials) on the spike history of the network.

The term:

$$X_k(\omega_0^{D-1}) = \frac{\theta - \mathcal{V}_k^{(det)}(\omega_0^{D-1})}{\sigma_k(\omega_0^{D-1})}, \quad (25)$$

contains the network spike history dependence of the neuron  $k$  at time  $D$ . More precisely, the term  $\mathcal{V}_k^{(det)}(\omega_0^{D-1})$  contains the deterministic part of the membrane potential of neuron  $k$  at time  $D$ , given the network spike history  $\omega_0^{D-1}$ , whereas  $\sigma_k(\omega_0^{D-1})$  characterizes the variance of the integrated noise in the neuron  $k$ 's membrane potential. We have:

$$\mathcal{V}_k^{(det)}(\omega_0^{D-1}) = \sum_{j=1}^N W_{kj} \eta_{kj}(\omega_0^{D-1}) + I_k \frac{1 - \gamma^{D-\tau_k(\omega_0^{D-1})}}{1 - \gamma}.$$

The first term is the network contribution to the neuron  $k$ 's membrane potential, where:

$$\eta_{kj}(\omega_0^{D-1}) = \sum_{l=\tau_k(\omega_0^{D-1})}^{D-1} \gamma^{D-1-l} \omega_j(l),$$

is the sum of spikes emitted by  $j$  in the past, with a weight  $\gamma^{D-1-l}$  corresponding to the leak decay of the spike influence as time goes on. The notation  $\tau_k(\omega_0^{D-1})$  means the last time before  $D-1$  where neuron  $k$  has fired, with the convention that this time is 0 if neuron  $k$  didn't fire between 0 and  $D-1$  in the block  $\omega_0^{D-1}$ . In the definition of  $\eta_{kj}(\omega_0^{D-1})$  we sum from  $\tau_k(\omega_0^{D-1})$ : this is because the membrane potential of neuron  $k$  is reset whenever  $k$  fires, hence loosing the memory of its past. Finally, in (25), we have:

$$\sigma_k^2(\omega_0^{D-1}) = \sigma_B^2 \frac{1 - \gamma^{2(D-\tau_k(\omega_0^{D-1}))}}{1 - \gamma^2}.$$

(see [7] for details)

## Explicit calculation of the canonical Maximum Entropy Potential

The goal now is to derive from (24) a canonical potential  $\mathcal{H}$  of the form (1) whose spike interactions terms  $h_i$ 's are functions of the network parameters: the synaptic weight matrix  $\mathcal{W}$  and the external stimulus  $\mathcal{I}$ ,  $h_i \equiv h_i(\mathcal{W}, \mathcal{I})$ .

Equation (14) gives a relation between the normalized potential and an equivalent non-normalized potential. From this equation, after considering the elimination of equivalent interactions is it possible to compute explicitly the values of the interaction terms  $h_i$ 's.

### Free energy:

From (15) and (24) we get the free energy:

$$-\phi(\omega^{(0)}) = \mathcal{P}[\mathcal{H}] = - \sum_{k=1}^N \log \left[ 1 - \pi \left( \frac{\theta - I_k \frac{1-\gamma^D}{1-\gamma}}{\sigma_B \sqrt{\frac{1-\gamma^{2D}}{1-\gamma^2}}} \right) \right].$$

### Local fields:

They are computed using equation (16). We consider the periodic orbit obtained by the  $R$ -circular shift of the block corresponding to the monomial  $\omega_i(D)$ . We have therefore to compute  $\phi(\omega^{(l_0)}) + \phi(\omega^{(l_1)}) + \dots + \phi(\omega^{(l_D)})$  using equation (24). To obtain this quantity we have to compute  $X_k$  (25) for all the blocks in the periodic orbit. Note that  $X_k$  does not depend on the last column of the blocks in the orbit. We abuse the notation by writing  $X_k(\omega^{(\sigma^n l)})$  instead of  $X_k(\omega_0^{(\sigma^n l)D-1})$ . The same holds for  $\eta_{kj}(\omega^{(\sigma^n l)})$  and  $\sigma_k(\omega^{(\sigma^n l)})$ . We obtain:

$$X_k(\omega^{(\sigma^n l)}) = \begin{cases} \frac{\theta - W_{ki} \gamma^{n-1} - I_k \frac{1-\gamma^D}{1-\gamma}}{\sigma_B \sqrt{\frac{1-\gamma^{2D}}{1-\gamma^2}}}, & 1 \leq n \leq R-1, k \neq i; \\ \frac{\theta - W_{kk} \gamma^{n-1} - I_k \frac{1-\gamma^n}{1-\gamma}}{\sigma_B \sqrt{\frac{1-\gamma^{2n}}{1-\gamma^2}}}, & 1 \leq n \leq R-1, k = i; \\ \frac{\theta - I_k \frac{1-\gamma^D}{1-\gamma}}{\sigma_B \sqrt{\frac{1-\gamma^{2D}}{1-\gamma^2}}}, & \forall k, n=0. \end{cases} \quad (26)$$

Combining equations (16), (24) and (26) we obtain:

$$h_i = \sum_{n=1}^{R-1} \sum_{k=1}^N \log \left[ 1 - \pi \left( X_k(\omega^{(\sigma^n l)}) \right) \right] + \sum_{k \neq i} \log \left[ 1 - \pi \left( X_k(\omega^{(\sigma^0 l)}) \right) \right] + \log \left[ \pi \left( X_i(\omega^{(\sigma^0 l)}) \right) \right] - R\phi(\omega^{(0)}). \quad (27)$$

which is an explicit function of synaptic weights and stimuli. Clearly:

- The "local field" of a neuron  $i$  depends non linearly on *all* stimuli (not only  $I_i$ ).
- It depends non linearly on the incoming synaptic weights connected to  $i$ .

### Pairwise interactions (instantaneous):

We get:

$$X_k(\omega^{(\sigma^{n_l})}) = \begin{cases} \frac{\theta - (W_{ki} + W_{kj}) \gamma^{n-1} - I_k \frac{1-\gamma^D}{1-\gamma}}{\sigma_B \sqrt{\frac{1-\gamma^{2D}}{1-\gamma^2}}}, & 1 \leq n \leq R-1, k \neq i, j; \\ \frac{\theta - (W_{kk} + W_{kj}) \gamma^{n-1} - I_k \frac{1-\gamma^n}{1-\gamma}}{\sigma_B \sqrt{\frac{1-\gamma^{2n}}{1-\gamma^2}}}, & 1 \leq n \leq R-1, k = i; \\ \frac{\theta - (W_{kk} + W_{ki}) \gamma^{n-1} - I_k \frac{1-\gamma^n}{1-\gamma}}{\sigma_B \sqrt{\frac{1-\gamma^{2n}}{1-\gamma^2}}}, & 1 \leq n \leq R-1, k = j; \\ \frac{\theta - I_k \frac{1-\gamma^D}{1-\gamma}}{\sigma_B \sqrt{\frac{1-\gamma^{2D}}{1-\gamma^2}}}, & \forall k, n=0. \end{cases} \quad (28)$$

Plugging (28) in (24) and using (18), one finally obtains  $J_{ij}$  as a explicit function of synaptic weights and stimulus.

Remarks:

- The “instantaneous pairwise” interaction  $J_{ij}$  depends not only on  $W_{ij}$ , but in all synaptic weights of neurons connected with  $i$  or  $j$ .
- It also depends in the stimulus of all neurons in the network.

### Pairwise interactions (1 time-step):

As mentioned in the previous section, in order to compute this term, the periodic orbit obtained by the  $R$ -circular shift of the block  $\omega^{(l_0)}$  corresponding to the monomial  $\omega_i(1)\omega_j(0)$  is not sufficient. We have therefore to use the periodic orbit obtained by our procedure (19). From  $\omega^{(l_1)}$  to  $\omega^{(l_4)}$  we have already computed their corresponding value  $X_k(\omega^{(\sigma^{n_l})})$  when computing the Local fields. From  $\omega^{(l_4)}$  to  $\omega^{(l_{R+2})}$  we just circularly shift  $\omega^{(l_4)}$ . We compute the corresponding  $X_k(\omega^{(\sigma^{n_l})})$ :

$$X_k(\omega^{(\sigma^{n_l})}) = \begin{cases} \frac{\theta - W_{ki} \gamma^{n-4} - W_{kj} \gamma^{n-3} - I_k \frac{1-\gamma^D}{1-\gamma}}{\sigma_B \sqrt{\frac{1-\gamma^{2D}}{1-\gamma^2}}}, & 4 \leq n \leq R+2, k \neq i, j; \\ \frac{\theta - W_{kk} \gamma^{n-4} - W_{kj} \gamma^{n-3} - I_k \frac{1-\gamma^{n-3}}{1-\gamma}}{\sigma_B \sqrt{\frac{1-\gamma^{2(n-3)}}{1-\gamma^2}}}, & 4 \leq n \leq R+2, k = i; \\ \frac{\theta - W_{kk} \gamma^{n-4} - W_{ki} \gamma^{n-3} - I_k \frac{1-\gamma^{n-2}}{1-\gamma}}{\sigma_B \sqrt{\frac{1-\gamma^{2(n-2)}}{1-\gamma^2}}}, & 4 \leq n \leq R+2, k = j; \end{cases} \quad (29)$$

We then apply equation (14) to obtain the desired term, from previously computed interaction terms.

A numerical illustration of our method is presented in figure (2). We start from the normalized potential (24) and construct the canonical equivalent potential. We then compare the conditional probability of patterns predicted by  $\mathcal{H}$  with the empirical probabilities inferred from a spike train generated by (24). This is just an

illustration, and not a systematic study. Note that this numerical analysis is limited to small  $N, R$  since the number of terms in  $\mathcal{H}$  grows exponentially fast, rendering intractable the method for  $NR \geq 20$ .

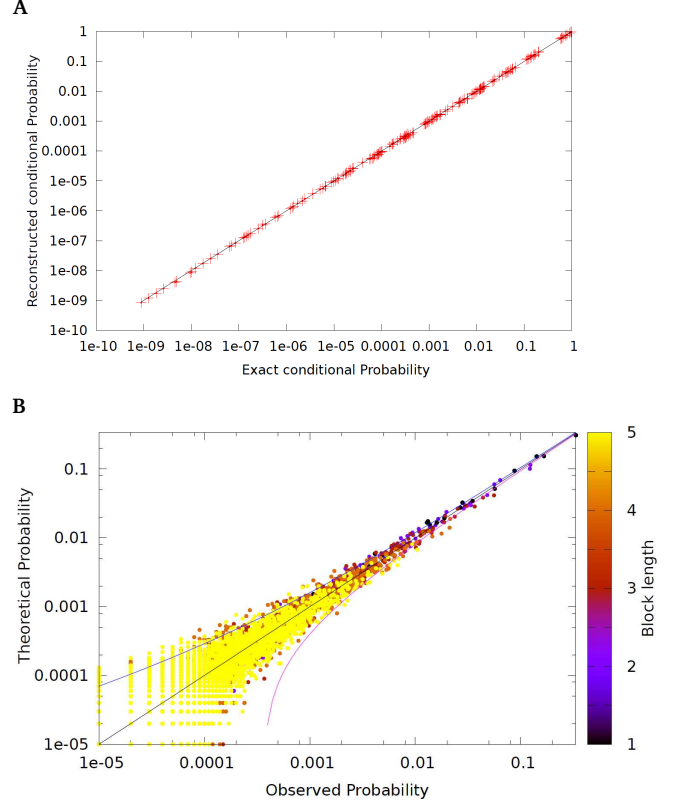


FIG. 2: (A). Exact conditional probabilities for blocks of range  $R$  obtained from the normalized potential (24), vs exact conditional probabilities associated with the potential (1). (B) Empirical probabilities of blocks  $\omega_0^k$ ,  $k = 1, \dots, 5$ , obtained from a discrete leaky integrate and fire spike train of size  $T = 10^5$  vs the probabilities of the same blocks predicted by the Gibbs distribution with potential  $\mathcal{H}(1)$ . Each dot stands for one of the  $2^{Nk}$  spatio-temporal patterns, where  $k$  is the block length. Diagonal shows equality. Confidence bounds (blue and red lines) correspond to fluctuations ruled by Central Limit Theorem. Plot is in log scale. This figure corresponds to  $N = 5, R = 3, \gamma = 0.2, \sigma_B = 0.2, \theta = 1, I_k = 0.7, k = 1, \dots, 5$ . The synaptic weights are random and sparse. Each neuron was randomly connected to other 2 neurons whose weights were drawn from a gaussian 0 mean and variance  $\frac{J^2}{N}$ . In this example  $J = 3$ .

### V. ANSWERING QUESTION 3 AND GENERAL CONCLUSION

When the normalized potential  $\phi$  is derived from a neuro-mimetic model (e.g. eq. (24)), it follows that the “local fields”  $\mathbf{h}_i$  depends non linearly on the complete stimulus  $\mathcal{I}$  (not only the stimulus applied to neuron  $i$ ), and the synaptic weights matrix  $\mathcal{W}$ . This is not that surprising. Even considering an Ising model of two

neurons with no memory, a strong favorable pairwise interaction between the two neurons will increase the average firing rate of both neurons, even in the absence of an external field. Likewise,  $J_{ij}$  depends on the whole synaptic weights matrix  $\mathcal{W}$  and not only on the connection between  $i$  and  $j$ . This example clearly shows that there is no straightforward relation between the so-called “functional connectivity” in Ising model  $J_{ij}$  and the neural synaptic connectivity ( $W_{ij}$ ).

As stated in the introduction a neuro-mimetic model with  $N$  neurons has  $O(N^2)$  parameters, whereas a MaxEnt model with  $N$  neurons and memory depth  $D$  has  $O(2^{NR})$  parameters  $h_l$  (canonical potential). Since the correspondence from neuro-mimetic models to a MaxEnt potential is exact we have two possibilities:

- (i) A large number of  $h_l$ ’s vanish.
- (ii)  $h_l$ ’s are related among them.

The two possibilities are actually not exclusive. Let us first address this question from the mathematical (dynamical systems) viewpoint which was the line followed up to now.

Consider a neuro-mimetic model with a well defined dynamics (e.g. (23)) and the associated normalized potential  $\phi = \phi(\mathcal{W}, \mathcal{I})$  (e.g. (24)). We may view a normalized potential as a point in a space: the coordinates of this point are fixed by  $\mathcal{W}, \mathcal{I}$ . A neuro-mimetic model corresponds therefore to the space of normalized potential of dimension  $O(N^2)$ . Using the same representation MaxEnt models with memory depth  $D$  span a space of dimension  $O(2^{NR})$ , but MaxEnt models equivalent to our neuro-mimetic model span a space of dimension  $O(N^2)$ . There is therefore a huge projection effect. Now eq (17), or, more generally, eq. (14) show that (ii) always functionally holds:  $h_l$  are non linear functions of  $\mathcal{W}, \mathcal{I}$  and are related to each others. This has a dramatic consequence. Assume that we want to fit (exactly) a neuro-mimetic model with a MaxEnt. We will need  $2^{NR}$  terms whereas  $O(N^2)$  are sufficient. This is exactly what happens in Fig 2. Now, we may have the hope that many  $h_l$ ’s are zero or close to zero. This is actually where MaxEnt models could make a breakthrough, showing that, in real spike trains many  $h_l$ ’s (almost) cancel would reveal a hidden law of nature. What is the hope for this? If we address this question from the dynamical systems viewpoint, there is no hope. Indeed in this context, one has to look for generic conditions for  $h_l$ ’s to vanish (case (i)). But it results from our analysis that the  $h_l$ ’s of a canonical potential corresponding to a neuro-mimetic model are *generically* non zero: considering e.g. *random* synaptic weights  $W_{ij}$ , the probability that some  $h_l$ ’s in (14) vanishes is indeed zero.

However, real neural networks are non generic: synaptic weights are not drawn at random but result from a long phylogenetic and ontogenetic evolution. When trying to “explain” spike statistics of real neural networks with the Maximum Entropy Principle, one is seeking some general laws that has to be expressed with relatively few phenomenological parameters in the potential (1). The hope is that many coefficients coming from real data are 0 or close to 0. This could explain the efficiency of pairwise MaxEnt models [4] for spike trains analysis (although this effect could also arise due e.g binning). Our method provides a way to detect this, if the l.h.s. in (14) is close to 0 [9].

Our method allows a mechanistic and causal understanding of the origin of correlations, in consequence, opens up new possibilities allowing a better understanding of the role of different neural network topologies and stimulus on the collective spike train statistics.

It is not limited to spike trains however and could also impact different areas of scientific knowledge where binary time series are considered.

## ACKNOWLEDGEMENTS

This work was supported by the French ministry of Research and University of Nice (EDSTIC), INRIA, ERC-NERVI number 227747, KEOPS ANR-CONICYT and European Union Project # FP7-269921 (BrainScales), Renvision # 600847. We thank the reviewers for constructive criticism.

- 
- [1] Y. Ahmadian, J. Pillow and L. Paninski *Neural Computation*, 23(1):46–96, 2011.
  - [2] V. I. Arnold. *Geometrical Methods in the Theory of Ordinary Differential Equations*. Springer-Verlag New York Inc., 1983.
  - [3] J. Besag *J. Roy Stat. Soc.*, 36(2):192–236, 1974.
  - [4] W. Bialek and R. Ranganathan. *arXiv:0712.4397*, 2007.
  - [5] R. Bowen. *Equilibrium States and the Ergodic Theory of Anosov Diffeomorphisms*. Springer-Verlag Berlin, 2008.
  - [6] D. R. Brillinger. *Biol Cybern*, 59(3):189–200, 1988.
  - [7] B. Cessac *J. Math. Biol.*, 62(6):863–900, 2011.
  - [8] B. Cessac and R. Cofré *J. Physiol. Paris*, 107(5):360–368, 2013.
  - [9] B. Cessac and R. Cofré. Research report, INRIA, 2013.
  - [10] R. Cofré and B. Cessac *Chaos Sol. and Frac.*, 50:13–31, 2013.
  - [11] E. J. Chichilnisky. *Network: Comput. Neural Syst.*, 12:199–213, 2001.
  - [12] S. Cocco, S. Leibler, and R. Monasson. *PNAS*, 106(33):14058–14062, 2009.
  - [13] R. Fernandez and G. Maillard. *J. Stat. Phys*, 118(3-4):555–588, 2005.
  - [14] E. Ganmor, R. Segev, and E. Schneidman. *The Journal of Neuroscience*, 31(8):3044–3054, 2011.

- [15] F. Gantmacher *The Theory of Matrices*. AMS Chelsea Publishing, 1959.
- [16] H.-O. Georgii. *Gibbs measures and phase transitions*. De Gruyter Studies in Mathematics:9. Berlin; New York, 1988.
- [17] W. Gerstner and W. Kistler. *Spiking Neuron Models*. Cambridge University Press, 2002.
- [18] E. Granot-Atedgi, G. Tkačik, R. Segev and E. Schneidman *Plos comp bio*, 9(3), 2013.
- [19] G. Grimmett *Bull Lon Math Soc*, 5(1), 1973.
- [20] J. M. Hammersley and P. Clifford. *unpublished*, 1971.
- [21] E.T. Jaynes. *Phys. Rev.*, 106:620, 1957.
- [22] G. Keller *Equilibrium States in Ergodic Theory*. Cambridge University Press, 1998.
- [23] A. Livšic. *Math. USSR- Izvestia*, (6):1278–1301, 1972.
- [24] O. Marre, S. El Boustani, Y. Frégnac, and A. Destexhe. *Phys. rev. Let.*, 102:138101, 2009.
- [25] M. Pollicott and H. Weiss. *Communications in Mathematical Physics*, 240:457–482, 2003.
- [26] F. Rieke, D. Warland, R. de Ruyter van Steveninck, and W. Bialek. *Spikes: Exploring the Neural Code*. Bradford Books, 1997.
- [27] M. Shadlen and W. Newsome *Curr. Opin. Neurobiol.*, 4:569–579, 1994.
- [28] E. Schneidman, M.J. Berry, R. Segev, and W. Bialek. *Nature*, 440(7087):1007–1012, 2006.
- [29] J. Shlens, G.D. Field, J.L. Gauthier, M.I. Grivich, D. Petrusca, A. Sher, A.M. Litke, and E.J. Chichilnisky. *Journal of Neuroscience*, 26(32):8254, 2006.
- [30] H. Soula and G. Beslon and O. Mazet *Neural Computation*, 18(1), 2006.
- [31] G. Tkačik, O. Marre, T. Mora, D. Amodei, M.J. Berry 2nd, and W. Bialek. *J Stat Mech*, page P03011, 2013.
- [32] J. C. Vasquez, O. Marre, A. G. Palacios, M. J Berry, and B. Cessac. *J. Physiol. Paris*, 106(3-4):120–127, 2012.
- [33] In the present context this property is ensured by the assumption  $\mathcal{H} > -\infty$  (sufficient condition).
- [34] When considering finite range potentials equilibrium states and Gibbs distributions are equivalent notions. This equivalence requires additional assumptions for infinite range potentials
- [35] Alternative constructions such as [24] has been proposed, but they require additional assumptions such as detailed balance
- [36] Thus, it corresponds to  $\gamma = 1 - \frac{dt}{RC}$  in the continuous-time LIF model.